

Dimensionality reduction

e.g. images \rightarrow input space $X \subseteq \mathbb{R}^{w \times h \times d}$ \rightarrow in principle we have such many degrees of freedom. This means that any value $\in \mathbb{R}^{w \times h \times d}$ could potentially be a valid element of the dataset. This is not true. There are often not so many degrees of freedom in the variability of the data in the dataset e.g. not all $w \times h \times d$ images contain real information, some are only noise, a bunch of random pixels.

We should make a distinction between the dimension of the input space and the actual variability of the data in the dataset.

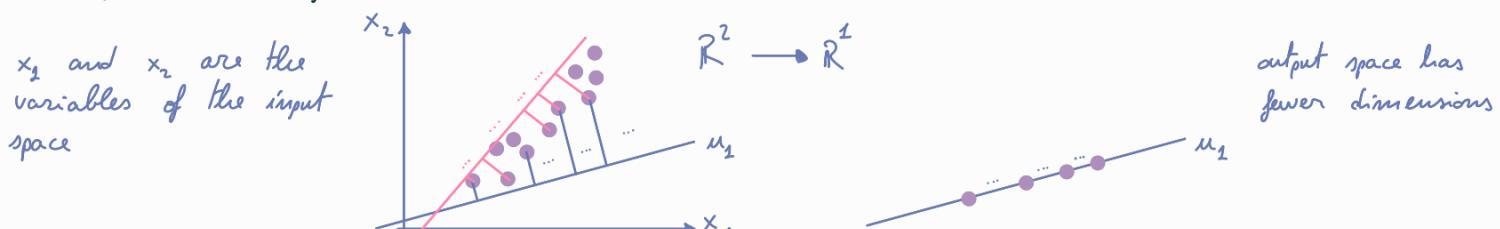
We want to study how to build and use a representation of the data in a much smaller dimensional space; the new representation should not lose important information. $\mathbb{R}^{w \times h \times d} \rightarrow \mathbb{R}^d \ll w \times h \times d$.

Dimensionality reduction: understand how to transform a problem in many dimensions to a problem in much less dimensions with the goal of keeping as much as possible the information. Sometimes information loss is inevitable.

Principal Component Analysis (PCA)

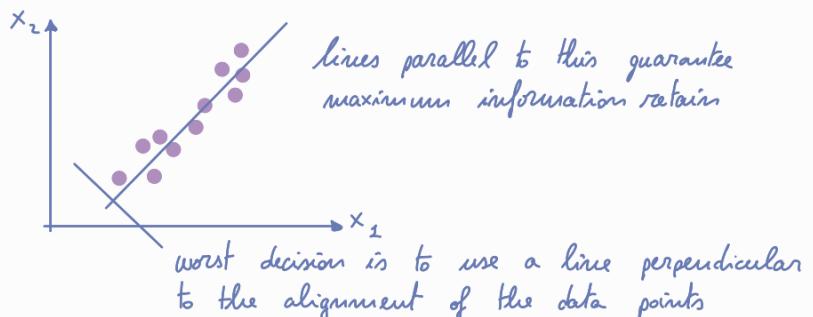
PCA is a dimensionality reduction method.

Given data $\{x_m\} \in \mathbb{R}^N$, the goal is to maximize variance after projection to some direction u . $\mathbb{R}^N \rightarrow \mathbb{R}^d$ with $d < N$. Output vector will be $\langle u_1^T x_m, \dots, u_d^T x_m \rangle$ with $\langle u_1, \dots, u_d \rangle$ are basis vectors for \mathbb{R}^d ($u_i^T u_j = 1 \forall i \neq j \rightarrow$ vectors are orthogonal)

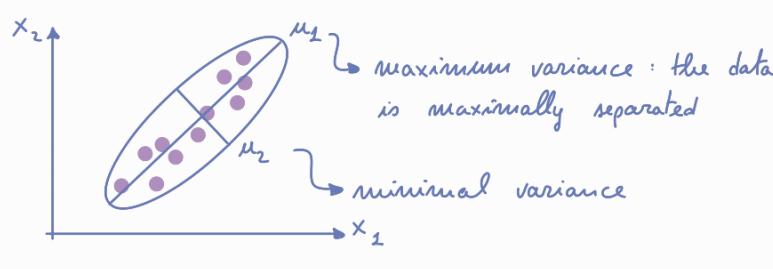


The point is: what is the best transformation? What transformation minimizes the error (maximum information retain)?

If the data is somewhat aligned (like in the previous example) the choice of the direction is easy

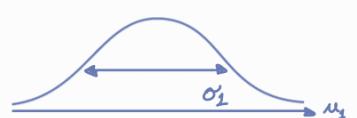


The general idea is to compute the covariance of the dataset and choose the direction of the maximum variance.



$\mathbb{R}^N \rightarrow \mathbb{R}^d$ components that maximize the variance

projection on u_1 :



projection on u_2 :



if you consider any other direction you will have an intermediate solution $\sigma_2 < \sigma < \sigma_1$

find the base vectors that yield the best transformation

$$\mathbb{R}^M \xrightarrow{\text{projection}} \mathbb{R}^D$$

$$\langle x_1, \dots, x_m \rangle \xrightarrow{\text{projection}} \langle u_1, \dots, u_D \rangle$$

we need to project from M dimensions to D dimensions

$$\langle u_1^T x_n, u_2^T x_n, \dots, u_D^T x_n \rangle \in \mathbb{R}^D$$

$$\langle u_1, \dots, u_D \rangle \text{ is a base} \rightarrow u_i^T u_j = 1 \quad \forall i, j, i \neq j$$

We just need to find $\langle u_1, \dots, u_D \rangle$, a particular set of base vectors that maximize the variance of the projected points. In order to do this we define:

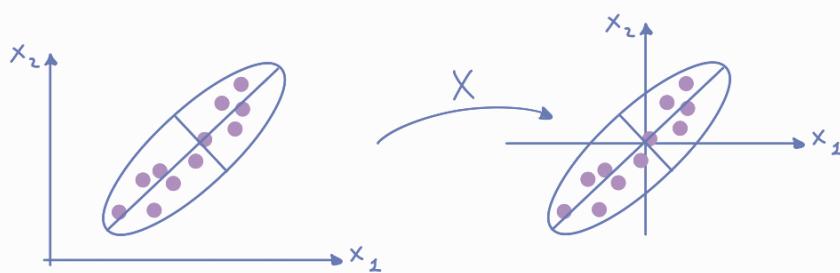
mean value of the data points

$$\bar{x} = \frac{1}{N} \sum_{m=1}^N x_m$$

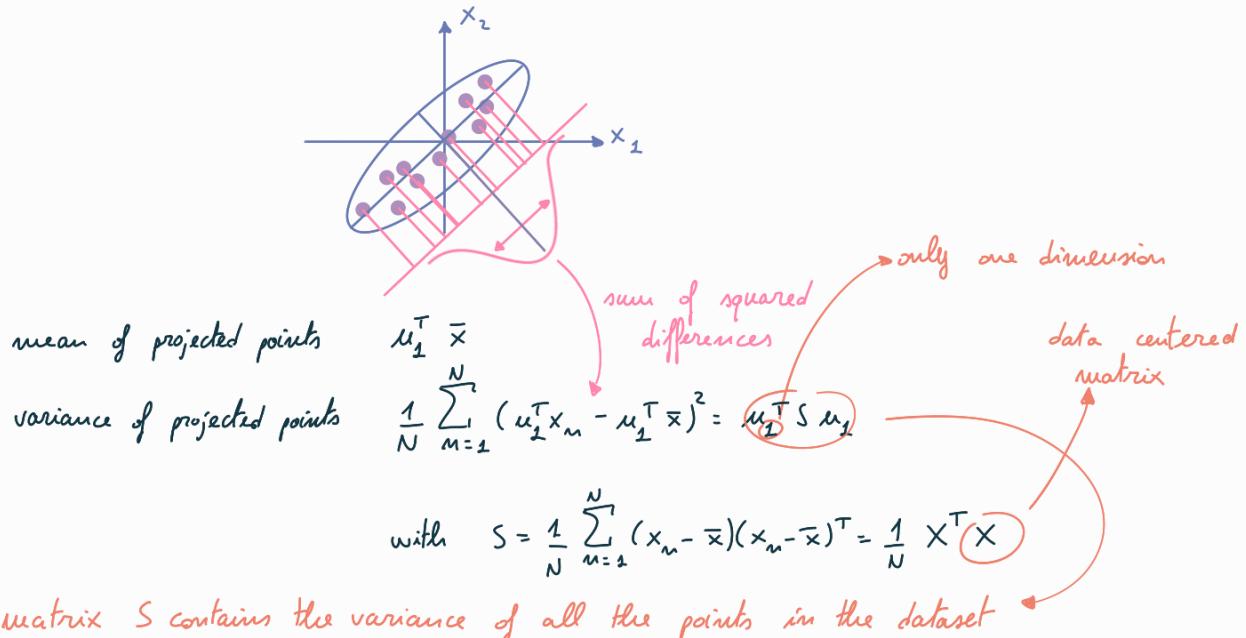
data-centered matrix

$$X = \begin{bmatrix} (x_1 - \bar{x})^T \\ \vdots \\ (x_i - \bar{x})^T \\ \vdots \\ (x_N - \bar{x})^T \end{bmatrix}$$

matrix with which we normalize the data so that the mean of the dataset matrix will be zero. This is equal to changing the reference system so that the mean is the new center.



Now compute the variance of the projected points; the new mean is the origin (thanks to the data centered matrix)



problem definition becomes: $\max_{u_1} u_1^T S u_1$

subject to constraint $u_1^T u_1 = 1$

e.g. solve with Lagrange multipliers

By solving this optimization problem we come to this solution:

$$u_1^T S u_1 = \lambda_1$$

λ_1 is an eigenvalue of S , u_1 is an eigenvector of $S \rightarrow$

when the variance of the projected points is maximum, we have this relation:

$$S u_1 = \lambda_1 u_1$$

maximum value (solution) corresponds to the eigenvalue of S , so we need to find the highest eigenvalue.

the maximum value $u_1^T S u_1$ is exactly λ_1 .

If we had a bidimensional space ($R^N \rightarrow R^2$), S will have two eigenvalues λ_1 and λ_2 and the solution of the maximum variance problem is the one that corresponds to the highest eigenvalue. This can be generalized in N dimensions.

If we order the eigenvalues in decreasing order, the first one will maximize the variance in the first dimension, the second one will maximize the variance in the second dimension and so on ...

Eigenvalues and Eigenvectors

For a square matrix A , an eigenvector and eigenvalue make this equation true :

$$A v = \lambda v$$

↑ eigenvalue
↑ eigenvector

$$\begin{matrix} \begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix} & \begin{bmatrix} 1 \\ 4 \end{bmatrix} & 6 \\ \text{matrix} & \text{eigenvector} & \text{eigenvalue} \end{matrix}$$
$$A \cdot v = \begin{bmatrix} -6 & 3 \\ 4 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 4 \end{bmatrix} = \begin{bmatrix} -6 \cdot 1 + 3 \cdot 4 \\ 4 \cdot 1 + 5 \cdot 4 \end{bmatrix} = \begin{bmatrix} 6 \\ 24 \end{bmatrix}$$
$$\lambda v = 6 \begin{bmatrix} 1 \\ 4 \end{bmatrix} = \begin{bmatrix} 6 \\ 24 \end{bmatrix}$$

Solution

$$u_2^T S u_2 = \lambda_2$$

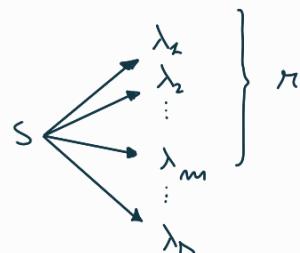
Variance is maximal when u_2 is the eigenvector corresponding to the largest eigenvalue λ_2 . This is called **principal component**. The direction u_2 is the corresponding eigenvector.

Approach :

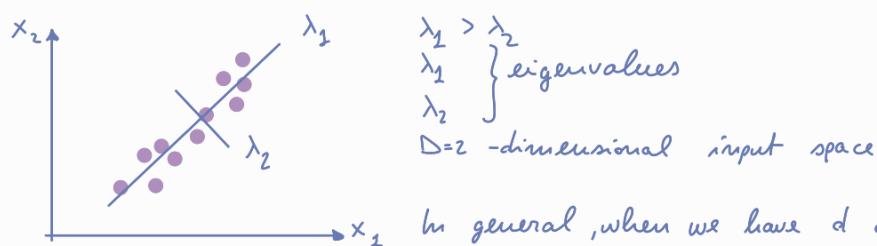
We start from the dataset of dimension D ; we want a dataset with dimension $M < D$.

- ① From the dataset we generate the data centered matrix
- ② From the new data we generate the mean and variance $\rightarrow u_2^T S u_2$ is the variance
- ③ We maximize the variance $\max_{u_2} u_2^T S u_2$

Solution to this maximization problem can be found computing the eigenvalues of S , ordering them from the highest to the lowest and taking the first M out of them.



This is still a **linear transformation**.



In general, when we have d dimensions with $d > 2$ we will have d eigenvalues. The highest eigenvalues correspond to the direction in which the projected points have maximum variance.

Once we project all the data on the first eigenvector we reduce the dataset from dimension d to dimension $d-1$; we can iterate this procedure. If you do this and repeat this procedure a second time you will obtain the second best eigenvector. In general, what we can do in order to reduce the space from D to M ($D > M$) is to :

- ① compute \bar{x} : mean of the data;
- ② compute S : covariance matrix of the data;
- ③ find M eigenvectors of S corresponding to the M largest eigenvalues.

Consider another formulation of the problem (transformation from a D -dimensional space to another D -dimensional space)

Given a d -dimensional dataset and a complete orthonormal d -dimensional basis such that dot product

$$u_i \cdot u_j = u_i^T u_j = \|u_i\| \|u_j\| \cos \theta$$

is $= 0$ when $\cos \theta = 0 \iff \theta = 90^\circ$

- $\iff u_i$ and u_j are perpendicular
- \iff linearly independent
- \iff part of a base for the space

each datapoint can be written as

representation of x in terms of the new base (u_1, \dots, u_d)



$$u_i^T u_j = \delta_{ij}$$

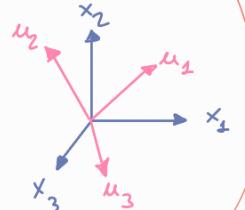
$$\delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

d -dimensional space + d -dimensional base:
we can consider the u vectors as a new reference space in which we want to transform the original dataset (new base)

Using the orthonormality property we have $\alpha_{mj} = x_m^T u_j$, hence:

$$x_m = \sum_{i=1}^D (\bar{x}_m^T u_i) u_i$$

The new basis has the same dimension of the original space.
This is only a transformation



Goal: approximate x_m using a lower dimensional representation

Split the above definition in two parts:

$$\tilde{x}_m = \sum_{i=1}^m z_{mi} u_i + \sum_{i=M+1}^D b_i u_i$$

first part is derived from the actual value in the dataset $\tilde{z}_{mj} = \bar{x}_m^T u_j \quad j=1, \dots, M$ \longrightarrow this part is correct

mean value $b_j = \bar{x}_m^T u_j \quad j=M+1, \dots, D$ \longrightarrow this part is an approximation

If we find a good way to approximate then we will only need the first M components.

\tilde{x}_m needs to be as close as possible to x_m .

Evaluate approximation error as:

$$J = \frac{1}{N} \sum_{m=1}^N \|x_m - \tilde{x}_m\|^2$$

The error in the approximation is given only by the last $D-M$ components, because the first M components are identical in x and \tilde{x} .

$$x - \tilde{x} = \sum_{M+1}^D [(x_m - \bar{x})^T u_i] u_i$$

the $D-M$ components are derived from the mean

We consider this $(x_m - \tilde{x}_m)$ as an error just for sample x_m

Putting together this procedure for each sample :

sum of squared errors for
all the samples in the dataset
where the error is considered
only on the last $D-r$ components

$$J = \left(\frac{1}{N} \sum_{m=1}^N \sum_{i=r+1}^D (x_m^T u_i - \bar{x}^T u_i)^2 \right) = \sum_{i=r+1}^D u_i^T S u_i$$

normalize for all the samples in the dataset

squared error covariance matrix

We have the same maximization problem that we had when we had to maximize the variance of the projection of each point, so we have the same solution :

$$S u_i = \lambda_i u_i$$

but since we have to minimize the error while before we had to maximize the variance, now we choose the lowest $D-r$ eigenvalues.

ordered from high to low $S = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_D \end{bmatrix}$

the first r are the ones that maximize the variance.
the last $D-r$ are the ones that minimize the approximation error.

We use the first r highest eigenvalues to project the dataset on this r dimensions and we automatically obtain that the result correspond to the minimization of the approximation error.